MICROCOPY RESOLUTION TEST CHART

NATIONAL BUREAU OF STANDARDS-1963-A

⑫ **LEVEL** *II*

⑪ 5 Dec 78

⑫ 14p.

⑥ CHI-SQUARE TESTS FOR THE
MULTINOMIAL DISTRIBUTION

⑩ KHURSHEED ALAM

DEPARTMENT OF MATHEMATICAL SCIENCES

CLEMSON UNIVERSITY

⑨ TECHNICAL REPORT 296

⑭ REPORT N104, TR-296

D D C

RECEIVED

JUN 6 1979

B

407 183

CHI-SQUARE TESTS FOR THE MULTINOMIAL DISTRIBUTION

Khursheed Alam

Clemson University

A simple proof of the asymptotic property of Chi-square tests, commonly used in the analysis of categorical data, is given for use as a note for instruction to first-year graduate students.

## 1.  Introduction

The Chi-square tests associated with the multinomial distribution are commonly used in the analysis of categorical data with reference to problems of specification, homogeneity of parallel samples, independence of attributes, etc. The asymptotic property of the tests, that is, the Chi-square distribution of the test statistics in large samples is generally known.  However, it has been our observation that many applied statisticians tacitly accept the asymptotic result without satisfying themselves with its proof.  It is also true that nearly all the text books in use on elementary and higher statistics either omit the proof or barely sketch it.  In this paper we outline a fairly simple proof of the fundamental result, for use as a note for the instruction to first-year graduate students and as a needed theory for the frequent application of Chi-square

tests by applied statisticians. The proof is essentially based on the contents of Chapters 5 and 6 of Rao (1966).
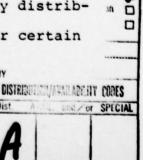
Consider a multinomial distribution $M(n,p)$ with $K$ cells, where $p = (p_1,\ldots,p_K)'$ denotes the vector of cell probabilities, $n = (n_1,\ldots,n_K)'$ denotes the vector of cell frequencies resulting from $n$ independent trials, $\sum_{i=1}^{K} p_i = 1$ and $\sum_{i=1}^{n} n_i = n$. A general problem in judging goodness of fit is to test whether the cell probabilities are specified functions of a fewer number of parameters whose values may be unknown. Let the cell probabilities be given functions $p_1(\theta),\ldots,p_K(\theta)$ of an unknown vector $\theta' = (\theta_1,\ldots,\theta_r)$, where $r < K$. To test the specification, it is a standard method to use the statistic

$$(1.1) \qquad T = \sum_{i=1}^{K} (n_i - np_i(\hat{\theta}))^2 / np_i(\hat{\theta})$$

where $\hat{\theta}$ is a consistent estimator of $\theta$, usually the maximum likelihood estimator. Next, assuming that the specification is true, consider the hypothesis that $\theta$ is given by $\theta_i = g_i(\alpha)$, where $g_1,\ldots,g_r$ are given functions, $\alpha' = (\alpha_1,\ldots,\alpha_s)$ and $s < r$. This hypothesis arises in a test of homogeneity of parallel samples and of independence in a contingency table. The statistic

$$(1.2) \qquad T^* = n \sum_{i=1}^{K} (p_i(\hat{\theta}) - p_i(\hat{\alpha}))^2 / p_i(\hat{\alpha})$$

is used to test the hypothesis, where $\hat{\alpha}$ denotes an estimate of $\alpha$ and $p_i(\alpha)$ denotes the value of $p_i(\theta)$ as a function of $\alpha$, under the given hypothesis. It is shown below that $T$ and $T^*$ are asymptotically distributed for large $n$ according to the Chi-square distribution under certain conditions.

## 2. Asymptotic Distribution of T and T*

First, consider the specification that the multinomial cell prob-
abilities are given functions $p_1(\theta),\ldots,p_K(\theta)$ of an unknown vector
$\theta = (\theta_1,\ldots,\theta_r)'$, where $r < K$. Let $\theta^0$ denote the true value of $\theta$. We
make the following assumptions:

(i) The functions $p_i(\theta)$ are continuous in $\theta$, admitting first order
partial derivatives which are continuous at $\theta^0$.

(ii) Given $\delta > 0$, there exists $\epsilon > 0$ such that $\displaystyle\inf_{|\theta-\theta^0|>\delta} N(\theta) > \epsilon$,

where
$$N(\theta) = \sum_{i=1}^{K} p_i(\theta^0) \log (p_i(\theta^0)/p_i(\theta)).$$

Let
$$M(\theta) = \sum_{i=1}^{K} \frac{n_i}{n} \log (p_i(\theta^0)/p_i(\theta)).$$

Consider the function $N(\theta)$ on the sphere $|\theta-\theta^0| = \delta$. Since $N(\theta)$ is contin-
uous in $\theta$, the infimum of $N(\theta)$ is attained on the sphere. Therefore, in
view of (ii), $N(\theta) \geq \epsilon$ for every point on the sphere. Since $\dfrac{n_i}{n}$ converges
in probability to $p_i(\theta^0)$ as $n \to \infty$, it follows that $M(\theta) > 0$ for all points
on the sphere with probability approaching 1 as $n \to \infty$ .

The log likelihood function is proportional to $\sum_{i=1}^{K} n_i \log p_i(\theta)$.
In view of (i) and the result given above, we have that for sufficiently
large n, the likelihood function has a local maximum inside the open
sphere $|\theta-\theta^0|<\delta$ at a point $\hat{\theta}$, say, which is a solution of the likelihood
equation

(2.1)
$$\sum_{i=1}^{K} \frac{n_i}{p_i(\theta)} \frac{\partial p_i(\theta)}{\partial \theta_j} = 0, \quad j = 1,\ldots,r.$$

Since $\delta$ can be made arbitrarily small, the maximum likelihood estimator
$\hat{\theta}$ is a consistent solution of the likelihood equation.

Let $I(\theta) = (M'(\theta))(M(\theta))$ denote the information matrix of the multinomial distribution, where

$$M(\theta) = (\frac{1}{\sqrt{p_i(\theta)}} \frac{\partial p_i(\theta)}{\partial \theta_j})$$

is a K x r matrix, and let $Z = M'V$, where $M = M(\theta^o)$ and

$$V' = (\frac{n_1 - np_1(\theta^o)}{\sqrt{np_1(\theta^o)}}), \ldots, \frac{n_K - np_K(\theta^o)}{\sqrt{np_K(\theta^o)}}).$$

By the central limit theorem, the asymptotic distribution of V is multivariate normal $N(\underset{\sim}{0}, I-\phi\phi')$, where $\underset{\sim}{0}$ denotes the null vector, I denotes the identity matrix and $\phi' = (\sqrt{p_1(\theta^o)}, \ldots, \sqrt{p_K(\theta^o)})$. The asymptotic distribution of Z is $N(0,I)$, where $I = I(\theta^o)$. Note that $M'\phi = \underset{\sim}{0}$.

Substituting $\hat{\theta}$ for $\theta$ in (2.1), the jth equation can be written as

$$(2.2) \qquad \sum_{i=1}^{K} \frac{n_i - np_i(\theta^o)}{\sqrt{n}p_i(\hat{\theta})} \frac{\partial p_i(\hat{\theta})}{\partial \hat{\theta}_j} = \sum_{i=1}^{K} \frac{\sqrt{n}(p_i(\hat{\theta}) - p_i(\hat{\theta}^o))}{p_i(\hat{\theta})} \frac{\partial p_i(\hat{\theta})}{\partial \hat{\theta}_j}.$$

In view of (i) we have

$$(2.3) \qquad p_i(\hat{\theta}) - p_i(\theta^o) = \sum_{\ell=1}^{r} (\hat{\theta}_\ell - \theta_\ell^o) \frac{\partial p_i(\theta^o)}{\partial \theta_\ell^o} + \eta|\hat{\theta} - \theta^o|$$

where $\eta \to 0$ as $\hat{\theta} \to \theta^o$. Since $\hat{\theta}$ is a consistent estimator of $\theta$, as shown above, the left side of (2.2) is asymptotically equivalent (converging in probability) to $Z_j$. Therefore, by the substitution of (2.3) on the right side of (2.2) we have that

$$z_j \overset{a}{\underset{\sim}{=}} \sum_{\ell=1}^{r} \sqrt{n}(\hat{\theta}_\ell - \theta_\ell^o) I_{\ell j}$$

where $\overset{a}{\underset{\sim}{=}}$ means "asymptotically equivalent to", and $I_{\ell j}$ denotes the $\ell j$th element of $I$. Hence

$$Z \overset{a}{\underset{\sim}{=}} \sqrt{n} \; I \; (\hat{\theta} - \theta^0)$$

or

$$(2.4) \qquad I^- Z \overset{a}{\underset{\sim}{=}} \sqrt{n} \; (\hat{\theta} - \theta^0)$$

where $I^-$ denotes a generalized inverse of $I$, given by $I \; I^- \; I = I$, and is equal to $I^{-1}$ if $I$ is non-singular.

Let A be a symmetric matrix with real elements, and let $X \overset{d}{\underset{\sim}{\to}} N (\mu, \Sigma)$, where $\overset{d}{\underset{\sim}{\to}}$ means "distributed as". If the covariance matrix $\Sigma$ is non-singular then it is known that the quadratic form $X' \; A \; X \overset{d}{\underset{\sim}{\to}} \chi^2_{\nu, \delta}$ (non-central Chi-square with $\nu$ degrees of freedom and non-centrality parameter $\delta$ ) if and only if $A \Sigma$ is idempotent, where $\nu = \text{Rank } A$ and $\delta = \frac{1}{2}\mu' A\mu$. If $\Sigma$ is singular, then the given condition is only sufficient and $\nu = \text{Rank } A\Sigma$ (see e.g. Graybill (1976), Theorem 4.7.1). If $A = \Sigma^-$ is a generalized inverse of $\Sigma$, given by $\Sigma \; \Sigma^- \Sigma = \Sigma$ , then $A \Sigma$ is idempotent and $\text{Rank } A \; \Sigma = \text{Rank } \Sigma$. Therefore, $X' \; \Sigma^- \; X \overset{d}{\underset{\sim}{\to}} \chi^2_{\nu, \delta}$ , where $\nu = \text{Rank } \Sigma$ and $\delta = \frac{1}{2} \mu' \; \Sigma^- \; \mu$.

Let $W = (W_1, \ldots, W_K)'$ where

$$W_i = \sqrt{n} \; (p_i(\hat{\theta}) - p_i (\theta^0)) / \sqrt{p_i(\theta^0)} \quad , \quad i = 1, \ldots, K.$$

From (2.3) we have that

$$(2.5) \qquad W \overset{a}{\underset{\sim}{=}} \sqrt{n} \; M \; (\hat{\theta} - \theta^0)$$

$$\overset{a}{\underset{\sim}{=}} M \; I^- \; Z \qquad \text{by (2.4)}$$

$$= M \; I^- \; M' \; V.$$

Note that $M \, I^- M'$ is idempotent. From (1.1) and (2.5) we have

$$T \overset{a}{\underset{\sim}{}} (V - W)' (V - W)$$

$$\overset{a}{\underset{\sim}{}} V' (I - M \, I^- M')V.$$

Now, V is asymptotically distributed as $N(\underset{\sim}{0}, I - \phi\phi')$, $(I - MI^-M')(I - \phi\phi')$ $= I - MI^-M' - \phi\phi'$ is idemptotent and

$$\text{Rank } (I - MI^-M' - \phi\phi) = \text{Trace } (I - MI^-M' - \phi\phi')$$
$$= K - 1 - \text{Rank } I = \beta, \text{say}.$$

Therefore, T is asymptotically distributed as $\chi_\beta^2$. If $I$ is of full rank then $\beta = K - r - 1$.

Next, consider the hypothesis that $\theta$ is given by $\theta_i = g_i(\alpha)$, $i = 1, \ldots, r$, as described in the previous section. Let $\nabla(\alpha) = (\partial\theta_i/\partial\alpha_j)$ denote the r x s matrix of the derivatives, and let $I^*(\alpha)$ denote the information matrix under the given hypothesis. Then

(2.6) $\qquad I^*(\alpha) = (\nabla(\alpha))' I(\theta) \nabla(\alpha)$

with $I(\theta)$ being expressed as a function of $\alpha$. Let $\alpha^0$ denote the true value of $\alpha$. Similarly, as in the preceding we have that

$$T^* \overset{a}{\underset{\sim}{}} V'(MI^-M' - M\nabla(\nabla'I\nabla)^-\nabla'M')V$$

where $\nabla = \nabla(\alpha^0)$. The matrix $(MI^-M' - M\nabla(\nabla'I\nabla)^-\nabla'M'( \; (I-\phi\phi') = (MI^-M' - M\nabla(\nabla'I\nabla)^-\nabla'M')$ is idempotent and

$$\text{Rank } (MI^-M' - M\nabla(\nabla'I\nabla)^-\nabla'M') = \text{Rank } I - \text{Rank } (\nabla'I\nabla)$$
$$= \gamma, \text{ say}.$$

Therefore, T* is asymptotically distributed as $\chi^2_\gamma$ . If $I$ and $\nabla$ are of full rank then $\gamma = r-s$.

We have shown that T and T* are asymptotically distributed according to the Chi-square distribution under the identifiability condition (i) and the continuity assumption (ii).

Remark: For the goodness of fit test where the cell probabilities are completely specified we have $\beta = K-1$. In this case $T \overset{a}{\sim} v'v \overset{d}{\sim} \chi^2_{K-1}$, asymptotically. For testing homogeneity of r parallel samples or independence of attributes in r x K contigency tables, we have $\gamma = (r-1)\ (K-1)$.

## References.

[1] Graybill, R. A. (1976). The Theory and Application of the Linear Model. Wadsworth Publishing Co., Belmont, California.

[2] Rao, C. R. (1966). Linear Statistical Inference and its Applications. Wiley Publications in Statistics.

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM | |
|---|---|---|---|
| 1. REPORT NUMBER<br>N104 | 2. GOVT ACCESSION NO. | 3. RECIPIENT'S CATALOG NUMBER | |
| 4. TITLE (and Subtitle)<br>Chi-square tests for the multinomial distribution | | 5. TYPE OF REPORT & PERIOD COVERED<br>Technical Report | |
| | | 6. PERFORMING ORG. REPORT NUMBER<br>296 | |
| 7. AUTHOR(s) | | 8. CONTRACT OR GRANT NUMBER(s)<br>N00014-75-C-0451 | |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br>Clemson University<br>Dept. of Mathematical Sciences<br>Clemson, South Carolina 29631 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS<br>NR 042-271 | |
| 11. CONTROLLING OFFICE NAME AND ADDRESS<br>Office of Naval Research<br>Code 436<br>Arlington, Va. 22217 | | 12. REPORT DATE<br>12-5-1978 | |
| | | 13. NUMBER OF PAGES<br>7 | |
| 14. MONITORING AGENCY NAME & ADDRESS(if different from Controlling Office) | | 15. SECURITY CLASS. (of this report)<br>Unclassified | |
| | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE | |

16. DISTRIBUTION STATEMENT (of this Report)

Approved for public release; distribution unlimited.

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)

18. SUPPLEMENTARY NOTES

19. KEY WORDS (Continue on reverse side if necessary and identify by block number)

Multinomial distributions; Chi-square test;
Contingency table.

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

This paper gives a proof of the asymptotic property of the Chi-square tests associated with the multinomial distribution, generally used in the analysis of categorical data, such as contingency tables.